# Hadoop WordCount Comparison

Kelompok 6
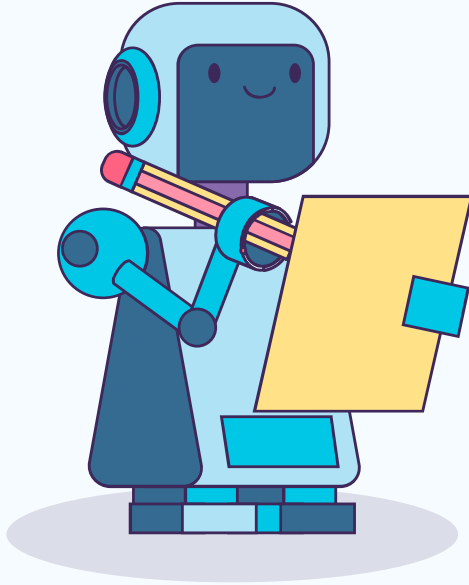
# Anggota

Muhammad Farrel Mirawan - 2106731554
Akmal Rabbani - 2106731610
Arka Brian Dewara - 2106731421
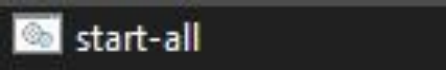Nevanda Fairuz Pahlevi - 2106731541

# 01

# Hadoop

start-all

← → C ⌂ ⓘ localhost:9870/dfshealth.html#tab-overview

Hadoop    Overview    Datanodes    Datanode Volume Failures    Snapshot    Startup Progress    Utilities ▾

# Overview 'localhost:9000' (active)

| Started: | Fri Jun 09 14:39:24 +0700 2023 |
|---|---|
| Version: | 3.2.2, r7a3bc90b05f257c8ace2f76d74264906f0f7a932 |
| Compiled: | Sun Jan 03 16:26:00 +0700 2021 by hexiaoqiao from branch-3.2.2 |
| Cluster ID: | CID-6f3bf16f-ee9a-4546-9a06-9fdc525ed978 |
| Block Pool ID: | BP-1393748948-192.168.43.1-1684981381356 |

← → C ⌂ ⓘ localhost:8088/cluster

# hadoop                                    # All Applications

▾ Cluster

**Cluster Metrics**

| | Apps Submitted | Apps Pending | Apps Running | Apps Completed | Containers Running | Used Resources | Total Resources | Reserved |
|---|---|---|---|---|---|---|---|---|
| | 0 | 0 | 0 | 0 | 0 | <memory:0, vCores:0> | <memory:8192, vCores:8> | <memory:0, |

About
Nodes
Node Labels
Applications
  NEW
  NEW_SAVING
  SUBMITTED
  ACCEPTED
  RUNNING
  FINISHED
  FAILED
  KILLED
Scheduler

**Cluster Nodes Metrics**

| Active Nodes | Decommissioning Nodes | Decommissioned Nodes | Lost Nodes | Unhealthy Nodes |
|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 |

**Scheduler Metrics**

| Scheduler Type | Scheduling Resource Type | Minimum Allocation | Maximum Allocation |
|---|---|---|---|
| Capacity Scheduler | [memory-mb (unit=Mi), vcores] | <memory:1024, vCores:1> | <memory:8192, vCores:4> |

▸ Tools

Show 20 entries

| ID | User | Name | Application Type | Queue | Application Priority | StartTime | LaunchTime | FinishTime | State | FinalStatus | Running Containers | Allocated CPU VCores | Allocated Memory MB | Allocated GPUs | Reserved CPU VCores | R |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | No data available in table | | | | | | | | |

Showing 0 to 0 of 0 entries

# WordCount Setup in Hadoop

```xml
<dependencies>
    <dependency>
        <groupId>org.apache.hadoop</groupId>
        <artifactId>hadoop-common</artifactId>
        <version>3.3.3</version>
    </dependency>
    <dependency>
        <groupId>org.apache.hadoop</groupId>
        <artifactId>hadoop-mapreduce-client-core</artifactId>
        <version>3.3.3</version>
    </dependency>
</dependencies>
```
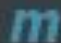
# WordCount Setup in Hadoop



```
Project ▾          ⊕ ≛ ⇟ ⚙ —      m pom.xml (WordCount) ×    © WordCount.class ×
WordCount D:\Games\WordCount          Decompiled .class file, bytecode version: 52.0 (Java 8)
  .idea                            19 ▶  public class WordCount {
  src                              20        public WordCount() {
    main                           21        }
      java                         22
        © WordCount                23 ▶      public static void main(String[] args) throws Exception {
      resources                    24            Configuration conf = new Configuration();
    test                           25            Job job = Job.getInstance(conf, jobName: "word count");
  target                           26            job.setJarByClass(WordCount.class);
  m pom.xml                        27            job.setMapperClass(TokenizerMapper.class);
External Libraries                 28            job.setCombinerClass(IntSumReducer.class);
Scratches and Consoles             29            job.setReducerClass(IntSumReducer.class);
                                   30            job.setOutputKeyClass(Text.class);
                                   31            job.setOutputValueClass(IntWritable.class);
                                   32            FileInputFormat.addInputPath(job, new Path(args[0]));
                                   33            FileOutputFormat.setOutputPath(job, new Path(args[1]));
```

# WordCount Setup in Hadoop

Create Jar using maven

| | |
|---|---|
| hadoop fs -mkdir /input_dir | Membuat folder input pada HDFS |
| hadoop fs -put /path/to/file.txt /input_dir | Meletakkan file text ke folder input HDFS |

# Menjalankan jar

```
PS D:\Games\WordCount> hadoop jar target/WordCount-1.0-SNAPSHOT.jar WordCount /input_dir /output_dir
```

# 02

# Hadoop Vs Java

# File yang akan diuji

| Name | Date modified | Type | Size |
|------|---------------|------|------|
| limaratusmb | 14/11/2014 5:05 | Text Document | 512.000 KB |
| satugb | 14/11/2014 5:05 | Text Document | 1.048.576 KB |
| satumb | 07/06/2023 23:58 | Text Document | 1.008 KB |
| sepuluhmb | 07/06/2023 23:59 | Text Document | 10.082 KB |
| seratusmb | 06/02/2016 19:38 | Text Document | 102.400 KB |

# Hadoop

# 1 Mb = 14.248 s

```
2023-06-08 00:39:00,545 INFO mapreduce.Job:  map 0% reduce 0%
2023-06-08 00:39:07,667 INFO mapreduce.Job:  map 100% reduce 0%
2023-06-08 00:39:14,785 INFO mapreduce.Job:  map 100% reduce 100%
2023-06-08 00:39:14,793 INFO mapreduce.Job: Job job_1686158454959_0003 completed successfully
```

# 10 Mb = 17.266 s

```
2023-06-08 00:35:19,777 INFO mapreduce.Job:  map 0% reduce 0%
2023-06-08 00:35:28,927 INFO mapreduce.Job:  map 100% reduce 0%
2023-06-08 00:35:36,021 INFO mapreduce.Job:  map 100% reduce 100%
2023-06-08 00:35:37,043 INFO mapreduce.Job: Job job_1686158454959_0002 completed successfully
```

# 100 Mb = 42.591 s

```
2023-06-08 00:41:32,217 INFO mapreduce.Job: Job job_1686158454959_0004 running in uber mode : false
2023-06-08 00:41:32,218 INFO mapreduce.Job:  map 0% reduce 0%
2023-06-08 00:41:49,453 INFO mapreduce.Job:  map 39% reduce 0%
2023-06-08 00:41:55,565 INFO mapreduce.Job:  map 51% reduce 0%
2023-06-08 00:42:02,654 INFO mapreduce.Job:  map 67% reduce 0%
2023-06-08 00:42:04,675 INFO mapreduce.Job:  map 100% reduce 0%
2023-06-08 00:42:14,808 INFO mapreduce.Job:  map 100% reduce 100%
```

# 500 Mb = 91.074 s

```
2023-06-08 00:52:58,328 INFO mapreduce.Job:  map 0% reduce 0%
2023-06-08 00:53:21,971 INFO mapreduce.Job:  map 6% reduce 0%
2023-06-08 00:53:23,082 INFO mapreduce.Job:  map 12% reduce 0%
2023-06-08 00:53:28,347 INFO mapreduce.Job:  map 17% reduce 0%
2023-06-08 00:53:29,398 INFO mapreduce.Job:  map 22% reduce 0%
2023-06-08 00:53:35,606 INFO mapreduce.Job:  map 23% reduce 0%
2023-06-08 00:53:37,695 INFO mapreduce.Job:  map 26% reduce 0%
2023-06-08 00:53:38,734 INFO mapreduce.Job:  map 31% reduce 0%
```

```
2023-06-08 00:54:11,953 INFO mapreduce.Job:  map 75% reduce 0%
2023-06-08 00:54:16,066 INFO mapreduce.Job:  map 77% reduce 0%
2023-06-08 00:54:17,088 INFO mapreduce.Job:  map 89% reduce 0%
2023-06-08 00:54:18,096 INFO mapreduce.Job:  map 95% reduce 0%
2023-06-08 00:54:19,108 INFO mapreduce.Job:  map 100% reduce 0%
2023-06-08 00:54:26,197 INFO mapreduce.Job:  map 100% reduce 100%
2023-06-08 00:54:29,254 INFO mapreduce.Job: Job job_1686158454959_0006 completed successfully
```

# 1 Gb = 191.408 s

```
2023-06-08 00:45:53,383 INFO mapreduce.Job:  map 0% reduce 0%
2023-06-08 00:46:24,480 INFO mapreduce.Job:  map 6% reduce 0%
2023-06-08 00:46:25,563 INFO mapreduce.Job:  map 7% reduce 0%
2023-06-08 00:46:29,795 INFO mapreduce.Job:  map 9% reduce 0%
2023-06-08 00:46:37,197 INFO mapreduce.Job:  map 10% reduce 0%
2023-06-08 00:46:42,451 INFO mapreduce.Job:  map 11% reduce 0%
2023-06-08 00:46:44,665 INFO mapreduce.Job:  map 17% reduce 0%
2023-06-08 00:46:48,964 INFO mapreduce.Job:  map 18% reduce 0%
2023-06-08 00:48:35,182 INFO mapreduce.Job:  map 85% reduce 25%
2023-06-08 00:48:41,329 INFO mapreduce.Job:  map 88% reduce 25%
2023-06-08 00:48:47,424 INFO mapreduce.Job:  map 90% reduce 25%
2023-06-08 00:48:53,516 INFO mapreduce.Job:  map 92% reduce 25%
2023-06-08 00:48:55,590 INFO mapreduce.Job:  map 96% reduce 25%
2023-06-08 00:48:56,595 INFO mapreduce.Job:  map 100% reduce 25%
2023-06-08 00:48:59,656 INFO mapreduce.Job:  map 100% reduce 64%
2023-06-08 00:49:03,687 INFO mapreduce.Job:  map 100% reduce 100%
2023-06-08 00:49:04,711 INFO mapreduce.Job: Job job_1686158454959_0005 completed successfully
```

# 10 Gb = 806 s

**Elapsed:** 13mins, 26sec

# Java

# 1 Mb = 110 ms

```
Run 1 - Runtime: 208 milliseconds

Run 2 - Runtime: 57 milliseconds

Run 3 - Runtime: 66 milliseconds

Average Runtime: 110 milliseconds
```

# 10 Mb = 509 ms

```
Run 1 - Runtime: 612 milliseconds

Run 2 - Runtime: 492 milliseconds

Run 3 - Runtime: 424 milliseconds

Average Runtime: 509 milliseconds
```

# 100 Mb = 5774 ms

```
Run 1 - Runtime: 5688 milliseconds

Run 2 - Runtime: 4385 milliseconds

Run 3 - Runtime: 7250 milliseconds

Average Runtime: 5774 milliseconds
```

# 500 Mb = 29594 ms

```
Run 1 - Runtime: 25921 milliseconds

Run 2 - Runtime: 39902 milliseconds

Run 3 - Runtime: 22960 milliseconds

Average Runtime: 29594 milliseconds
```

# 1 Gb = 51647 ms

```
Run 1 - Runtime: 58386 milliseconds

Run 2 - Runtime: 48680 milliseconds

Run 3 - Runtime: 47877 milliseconds

Average Runtime: 51647 milliseconds
```
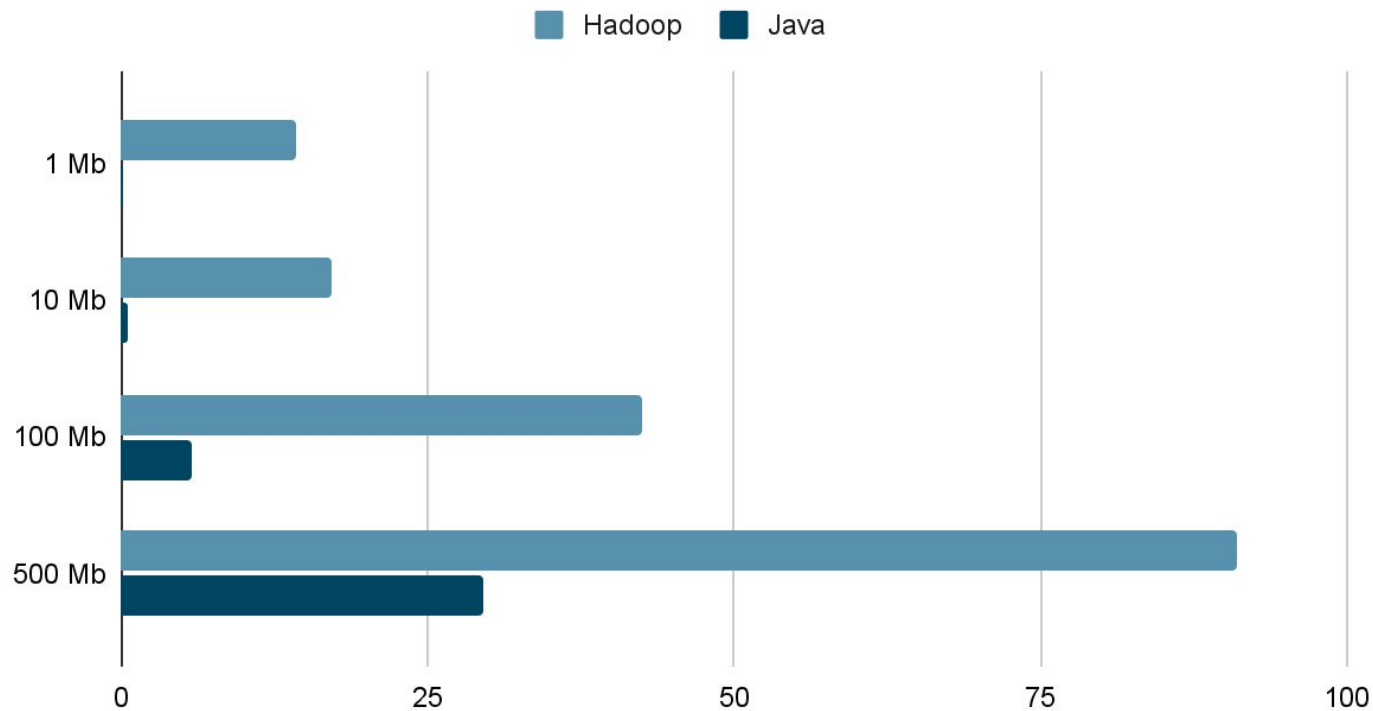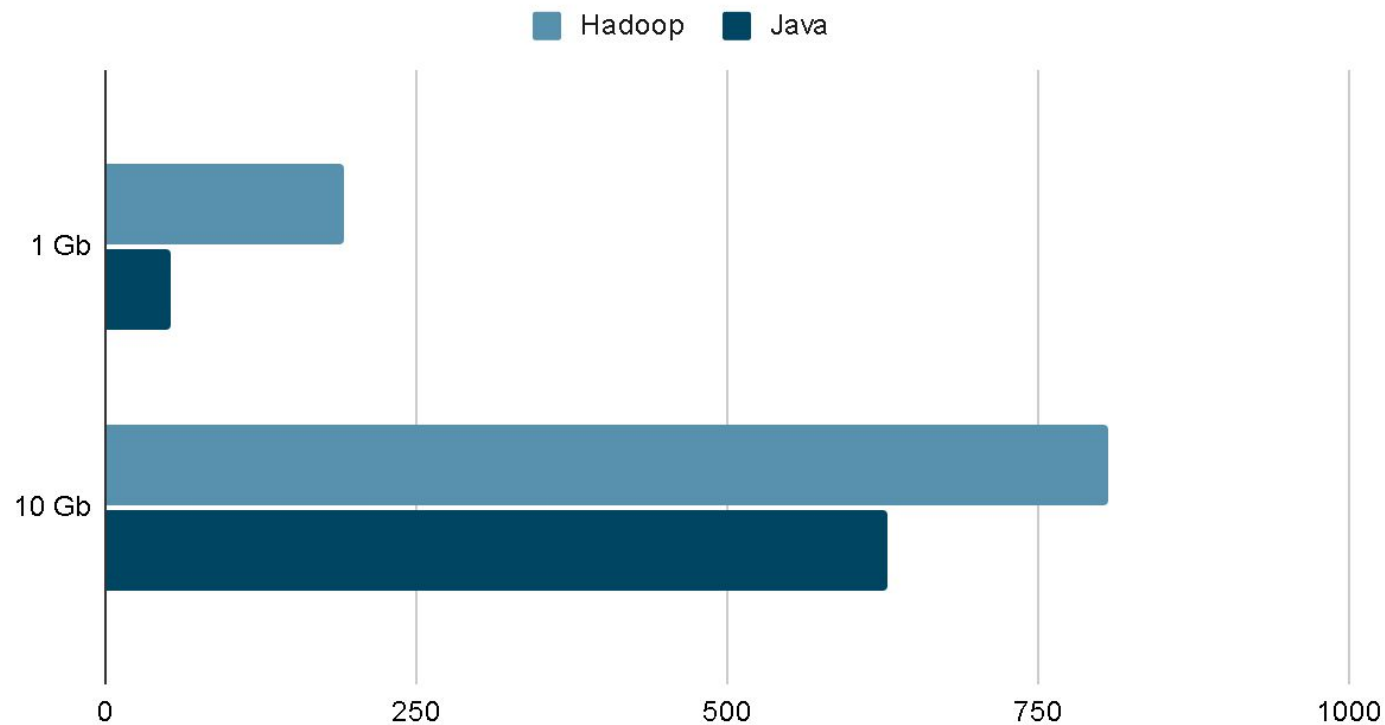
# 10 Gb = 628307 ms

```
Run 1 - Runtime: 628307 milliseconds

Average Runtime: 628307 milliseconds
```

# Hadoop vs Java



Legend: Hadoop, Java

| | |
|---|---|
| 1 Mb | |
| 10 Mb | |
| 100 Mb | |
| 500 Mb | |

Axis: 0, 25, 50, 75, 100

Hadoop vs Java (1Gb and 10 Gb)

Legend: Hadoop, Java

1 Gb, 10 Gb

X-axis: 0, 250, 500, 750, 1000

# Analisis

Berdasarkan hasil runtime word count yang diperoleh dari file dengan ukuran 1Mb, 10Mb, 100Mb, 500Mb, 1Gb, dan 10 Gb, dapat dilihat bahwa untuk melakukan word count dengan menggunakan hadoop memerlukan waktu yang cenderung lebih banyak dibandingkan dengan yang tidak menggunakan hadoop, dengan bahasa pemrograman Java. Hal ini dapat dilihat dari grafik yang dihasilkan, yang mana runtime untuk melakukan word count dengan hadoop selalu lebih besar. Jadi, file size yang diuji tidak menentukan kecepatan dari runtime yang dilakukan untuk word count tersebut. Hal ini disebabkan karena Hadoop memiliki overhead yang menjadikannya tidak cocok untuk file size kecil melainkan hadoop ini jauh lebih cocok digunakan apabila data yang diatur nya memiliki file size yang sangat besar.

**THANK YOU**